

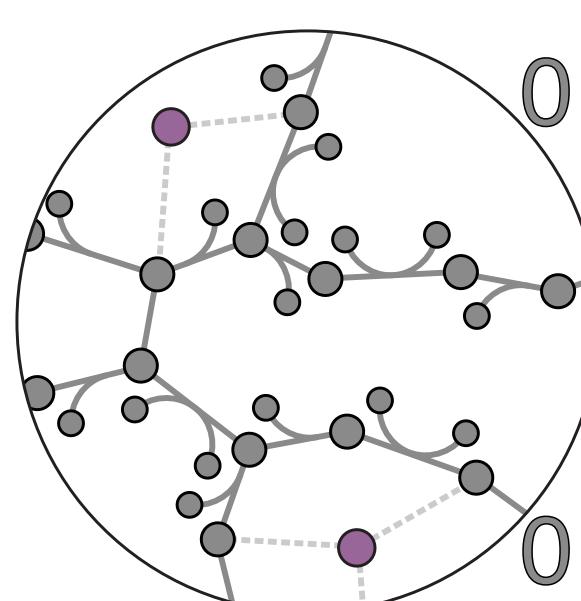


ALE Analytics - A microbial data warehouse and analytical suite

Patrick V. Phaneuf¹, Dennis Gosting², Bernhard O. Palsson^{2,3}, Adam M. Feist^{2,3}

1) Department of Computer Science and Engineering, University of California at San Diego, La Jolla 92093-0404
2) Department of Bioengineering, University of California San Diego, La Jolla, CA, 92093, United States of America,
3) Novo Nordisk Foundation Center for Biosustainability, Technical University of Denmark, 2800, Lyngby, Denmark

systems biology research group

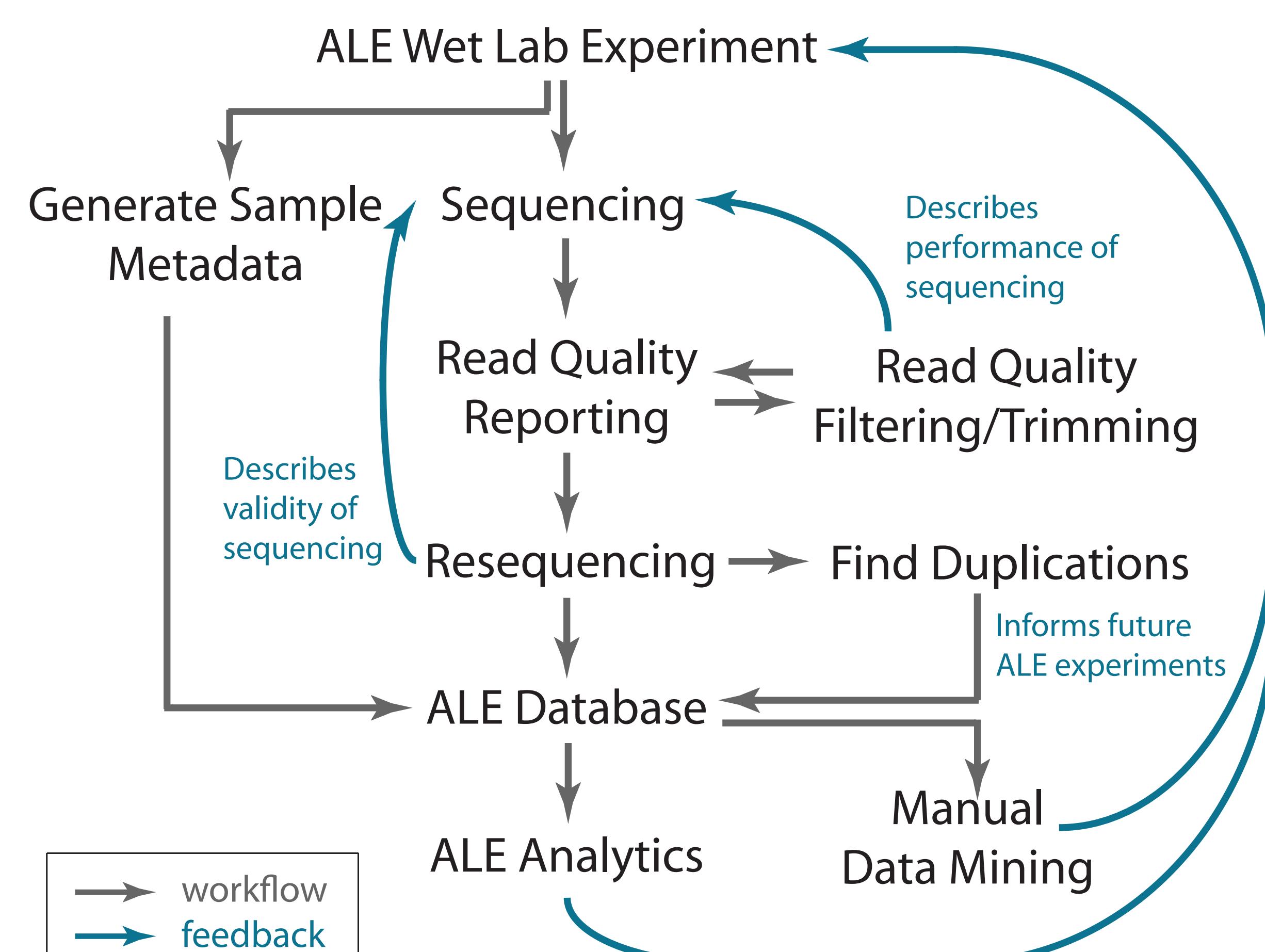


Abstract

Adaptive Laboratory Evolution (ALE) is a powerful experimental tool used by academic and industrial labs for a number of applications. ALE experiments generate a substantial amount of data, coming in the form of sequencing reads, alignment reports, and sample metadata. As ALE experiments scale to include more samples, the task of managing this data comes at higher costs due to the effort necessary to organize and integrate data into a format that describes the evolution process succinctly. In this work, we describe the development, deployment and iteration of an 'ALE Analytics' pipeline and web platform that streamlines the necessary ALE experiment data post-processing, manages experiment data, and produces interactive reports that detail an ALE experiment. Our design has been primarily driven by the need to consolidate large amounts of ALE experimental data in such a way to describe the quality of the sample sequencing, adaptive mutations in evolved strains, the context of mutations via their metadata (i.e., culturing environments, strain properties), and related mutations found in other experiments housed in the database. We have done so by leveraging a full stack of technologies that enable the parsing and databasing of experiment data, the execution of automated analysis on said data and the generation of web accessible reports. Future efforts will take full advantage of this developing platform to enable more depth and breadth of ALE experiment analysis with quicker turnaround.

ALE Mutation Processing Pipeline

Our ALE pipeline leverages industry established tools to process and refine ALE replicate sequencing data to identify genomic perturbations and integrate them into a database that can be mined for mutational trends.



Features and Advantages

- Supports multiple bacterial strains.
- Leverages established tools: Breseq, FastQC, FASTX-Toolkit, etc.
- Automated multi-sample processing for specific stages.
- Stage dependent feedback.
- Feedback provides opportunity for experiment and data refinement.

Contact

pphaneuf@eng.ucsd.edu
afeist@ucsd.edu
systemsbiology.ucsd.edu

References

- LaCroix RA et al. Appl Environ Microbiol 2014, 81:17-30
- Guzman GL et al. PNAS 2015, 112:929-934
- DOI 10.5281/zenodo.20980
- DOI 10.5281/zenodo.14561

ALE Analytics Web Platform Features

The ALE Analytics web platform integrates our ALE mutation analysis with published multi-omics data and better enables identification of causal mutations in ALE experiments.

ALE Database

Database containing the mutational data and metadata generated by the ALE data processing pipeline. This database is the basis for all analytical features implemented within the ALE Analytics web platform

ALE Analytics

A web interface reporting ALE experiment data and trends found through automated datamining.

Replicate Metadata

Contextualizing mutations using descriptions of each ALE sample: strain, genetic perturbations, media, substrate, temperature, etc.

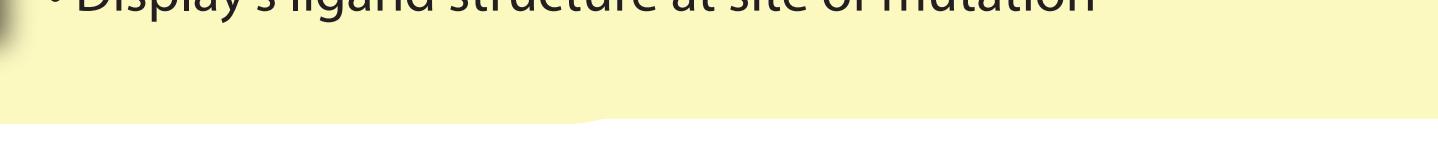
ALE Database trend reporting

Mutation Search

Search the ALE Database by mutation specific features.

Genes

Gene specific information:
 • All ALE Database mutations of a particular gene.
 • 3D rendering of a gene's product.
 • Highlight structure of a gene's product effected by a mutation.
 • Display ligand structure at site of mutation



Future Directions

- Use whole ALE Database statistics to rank the novelty of causal mutations. An ALE experiment's mutated gene found to have few mutations within the ALE Database can be considered relatively stable and is therefore more significant than a mutated gene with many ALE Database mutations.
- Datamine for novel regulatory genomic regions.

- Datamine for pairwise mutations shared among different ALE Experiments and establish their relationship.
- Integrate ALE experiment fitness data to provide additional context for identifying causal mutations.

Common Mutations

A report of ALE samples and their sets of shared mutations. Used to find samples most mutationally similar to a currently selected sample.

Frequently Mutated Genes

A report of an ALE experiment's multiple replicates and their sets of mutations that:

- Affect the same gene within a replicate.
 - Affect the same gene among parallel replicates.
- Highlights which genes are more frequently mutated within an ALE experiment, potentially elucidating the functional categories that were more causal.

Mutation Fixation

A report of ALE samples and their mutations that emerged and fix within an ALE. Presents which mutations were kept by the dominant populations between flasks and therefore may have lead to fitness benefits.

Dashboard

General statistics on entire ALE Database including:

- Most frequently mutated genes
- Most frequent mutations
- Mutation type counts (SNP, INDEL, etc.)

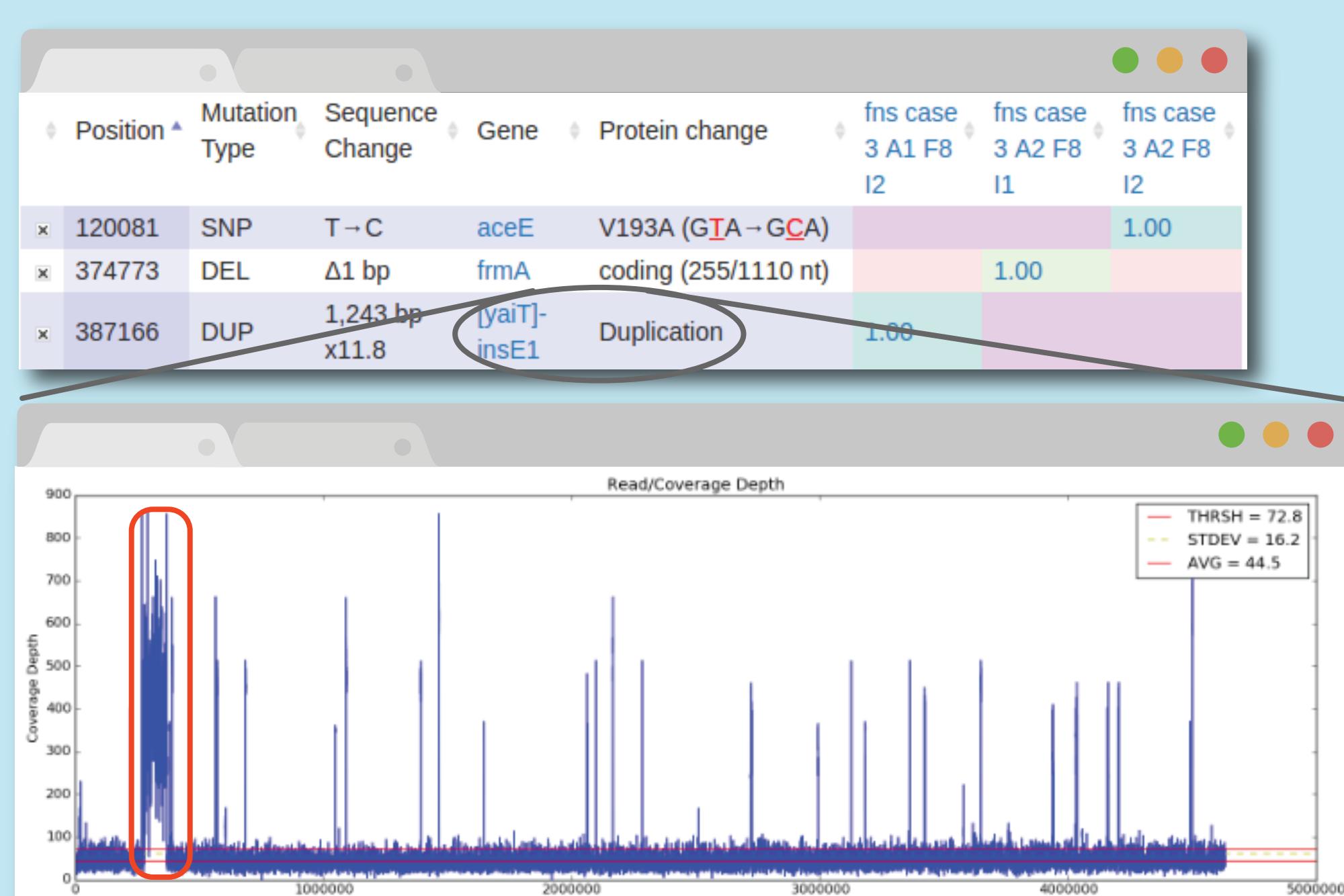
ALE Analytics Unique Mutation Stats		ALE Analytics Unique Functional Change Stats	
Mutation Type	Count	Functional Change Type	Count
Single Base Substitutions	11163	Intergenic	3564
Multiple Base Substitutions	41	Noncoding	351
Deletions	959	Pseudogene	337
Insertions	377	Synonymous	2044
Mobile Element Insertions	137	Nonsynonymous	5090
Amplifications	3		
Gene Conversions	0		
Inversions	0		
Duplications	1593		

Top Genes

Gene	Count
rmlH	2521
rrsH	1652
rrlA	1144
carB	1055

Duplications

Genomic duplication and amplification finding based on resequencing alignment read-depth. Can find large areas of duplication that are not otherwise found with current resequencing tools.



Experiment Statistics

Statistics and visualizations that can lend intuition to an ALE experiment's sequencing data quality, mutational "hotspots", most frequently mutated genes and most frequent mutations.

